

MASSACHUSETTS INSTITUTE OF TECHNOLOGY  
ARTIFICIAL INTELLIGENCE LABORATORY

AIM 355

March 1976

**ARTIFICIAL INTELLIGENCE -- a personal view**

by

**D. Marr**

**ABSTRACT.** The goal of A. I. is to identify and solve useful information processing problems. In so doing, two types of theory arise. Here, they are labelled Types 1 and 2, and their characteristics are outlined. This discussion creates a more than usually rigorous perspective of the subject, from which past work and future prospects are briefly reviewed.

*ARTIFICIAL INTELLIGENCE -- a personal view*

Artificial Intelligence is the study of complex information processing problems that have their roots in some aspect of biological information processing. The goal of the subject is to identify useful information processing problems, and give an abstract account of how to solve them. Such an account is called a *method*, and it corresponds to a theorem in mathematics. Once a method has been discovered for solving a problem, the final stage is to develop algorithms that implement the method. Choice of an algorithm for a given method usually depends upon the hardware in which the process is to run, and there may be many algorithms that implement the same method. A method on the other hand characterizes the solution to a problem at a more abstract level. The term "method" is in fact taken from Jardine & Sibson's (1971) work on cluster analysis, where they decomposed the subject in precisely this way.

A *result* in Artificial Intelligence thus consists of the isolation of a particular information processing problem, and the statement of a method for solving it. Some judgement has to be applied when deciding what constitutes a method - something based on exhaustive search for chess would clearly not qualify. The kind of judgement that is needed seems to be rather similar to that which decides whether a result in mathematics amounts to a substantial new theorem, and I do not feel uncomfortable about having to leave the basis of such judgements unspecified. The important point is that once a method has been established for a particular problem it never has to be done again, and in this respect a result in A.I. behaves like a result in mathematics or any of the hard natural sciences. It is conceivable that two different methods may exist for the same problem, just as a mathematical theorem may admit of two completely different proofs, but this circumstance has not yet occurred. New algorithms for implementing a known method may subsequently be devised without throwing substantial new light upon the method, just as discovering the fast fourier transform algorithm (Cooley & Tukey 1965) shed no new light on the nature of fourier analysis.

This view of what constitutes a result in A.I. is probably acceptable to most scientists. Chomsky's (1965) notion of a "competence" theory for English syntax is precisely what I mean by a "method" for that problem. Both have the quality of being little concerned with the gory details of algorithms that must be run to express the competence (i.e. to implement the method). That is not to say that devising suitable algorithms will be easy, but it is to say that before one can devise them, one has to know what exactly it is that they are supposed to be doing, and this information is captured by the method. When a problem decomposes in this way, I shall refer to it as having a *Type 1* theory.

The fly in the ointment is that while many problems of biological information processing have a Type 1 theory, there is no reason why they should all have. This can happen when a problem is solved by the simultaneous action of a considerable number of processes, *whose interaction is its own simplest description*, and I shall refer to such a situation as a *Type 2* theory. One promising candidate for a Type 2 theory is the problem of predicting how a protein will fold. A large number of influences act on a large polypeptide chain as it flaps and flails in a medium. At each moment only a few of the possible

interactions will be important, but the importance of those few is decisive. Attempts to construct a simplified theory must ignore some interactions; but if most interactions are crucial at some stage during the folding, a simplified theory will prove inadequate. [This situation is similar to the situation in mathematics, where any system that is complex enough to model arithmetic must essentially contain it]. Interestingly, the most promising studies of protein folding are currently those that take a brute force approach, setting up a rather detailed model of the amino acids, the geometry associated with their sequence, hydrophobic interactions with the circumambient fluid, random thermal perturbations *etc.*, and letting the whole process run until a stable configuration is achieved (Levitt & Warshel 1975).

The principal difficulty in A.I. is that one can never be quite sure whether a problem has a Type 1 theory. If one is found, well and good; but failure to find one does not mean that it does not exist. Most A.I. programs have hitherto amounted to Type 2 theories, and the danger with such theories is that they can bury crucial decisions, that in the end provide the key to the correct Type 1 theory, beneath the mound of small administrative decisions that are inevitable whenever a concrete program is designed. This phenomenon makes research in A.I. difficult to pursue and difficult to judge. If one shows that a given information processing problem is solved by a particular, neatly circumscribed method, then that is a secure result. If on the other hand one produces a large and clumsy set of processes that solves a problem, one cannot always be sure that there isn't a simple underlying method for one or more related problems whose formulation has somehow been lost in the fog. With any candidate for a Type 2 theory, much greater importance is attached to the performance of the program. Since its only possible virtue might be that it works, it is interesting only if it does. Often, a piece of A. I. research has resulted in a large program without much of a theory, which commits it to a Type 2 result, but that program either performs too poorly to be impressive or (worse still) has not even been implemented. Such pieces of research have to be judged very harshly, because their lasting contribution is negligible.

Thus we see that as A.I. pursues its study of information processing problems, two types of solution are liable to emerge. In one, there is a clean underlying theory in the traditional sense. Examples of this from vision are Horn's (1970) method for obtaining shape from shading, the notion of the Primal Sketch (Marr 1975), Ullman's (1975) method for detecting light sources, T. O. Binford's generalised cylinder representation, on which Marr & Nishihara's (1975) theory of the internal representation and manipulation of 3-D structures was based, and various other pieces of work currently in progress at our laboratory. The characteristic of these results is that they often lie at a relatively low level in the overall canvas of intellectual functions, a level often dismissed with contempt by those who purport to study "higher, more central" problems of intelligence. Our reply to such criticism is that low-level problems probably do represent the easier kind, but that is precisely the reason for studying them first. When we have solved a few more, the questions that arise in studying the deeper ones will be clearer to us.

But even relatively clean Type 1 theories such as these involve Type 2 theories as well. For example, Marr & Nishihara's 3-D representation theory asserts that the

deep underlying structure is essentially that of a stick figure, and that this representation is explicitly manipulated during the analysis of an image. Such a theory would be little more than speculation unless it could also be shown that such a description may be computed from an image and can be manipulated in the required way. To do so involves several intermediate theories, for some of which there is hope of eventual Type 1 status, but others look intractably of Type 2. For example, it now looks as though a Type 1 theory can be constructed for the problem of segmenting a bounding contour into units that correspond to the three-dimensional object's generalised cylinder components; but we see little prospect of deriving a Type 1 theory for the basic grouping processes that operate on the primal sketch to help separate figure from ground. A number of central ideas are involved there, but much of the difficulty resides in the engineering details, and we see no reason why those could or should be derivable from some single underlying theory; they amount to a procedural representation of many facts about images that derive ultimately *via* evolution from the cohesion and continuity of matter in the physical world. Many kinds of knowledge and different techniques are involved; one just has to sort them out one by one. As each is added the performance of the whole improves, and the complexity of the images that can be handled increases.

We have already seen that to search for a Type 2 theory for a problem may be dangerous if in fact it has a Type 1 theory. This danger is most acute in premature assaults on a high-level problem, for which few or none of the concepts that underlie its eventual Type 1 theory have yet been developed, and the consequence is a complete failure to formulate correctly the problems that are in fact involved. But it is equally important to realise that the opposite danger exists lower down. For example, in our current theory of visual processing, the notion of the primal sketch seems respectable enough, but one might have doubts about the aesthetics of the grouping processes that decode it. There are many of them, their details are somewhat messy; and seemingly arbitrary preferences occur (e.g. for vertical or horizontal organizations). A clear example of a Type 2 theory is our assertion that texture-vision discriminations rest on these grouping processes and first-order discriminations applied to the primal sketch of the image (Marr 1975). As such, it is less attractive than Julesz's (1975) clean (Type 1) theory that textured regions are discriminable if and only if there is a difference in the first or second-order statistics of their intensity arrays. In practice the two theories overlap considerably, because many of the processes we use are so simple that they can be represented as second-order operations on the image (via their Volterra series expansion, for example). The two theories are not identical however, and at present ours accounts for the counter-examples that Julesz found to his own theory.

I feel that one should not be too distressed if the messier theory turns out to be correct. We already know that separate modules must exist for computing other aspects of visual information, -- motion, stereoscopy, fluorescence, color, - and there is no reason why they should all be based on a single theory. Indeed one would *a priori* expect the opposite; as evolution progressed, new modules come into existence that can cope with yet more aspects of the data, and as a result kept the animal alive in ever more widely ranging circumstances. The only important constraint is that the system as a whole should be roughly

modular, so that new facilities can be added easily.

So, especially at the more peripheral stages of sensory information processing, and perhaps also more centrally, one should not necessarily give up if one fails to find a Type 1 theory -- there may not be one. More importantly even if there were, there would be no reason why that theory should bear much relation to the theory of more central phenomena. In vision for example, the theory that says 3-D representations are based on stick figures and shows how to manipulate them, is independent of the theory of the primal sketch, or for that matter of most other stages en route from the image to that representation. In particular, it is especially dangerous to suppose that an approximate theory of a peripheral process has any significance for higher level operations. For example, because Julesz's second-order statistics idea is so clean and so neatly fits much data, one might be tempted to ask whether the idea of second order interactions is in some way central to higher processes. In doing so one should bear in mind that the true explanation of texture visual discrimination may be quite different in nature even if the theory is very often a correct predictor of performance.

The reason for drawing this point out at such length is that it bears upon another issue, namely the type of theory that the grammar of natural language might have. The purpose of human language is presumably to transform a datastructure that is not inherently one-dimensional into one-dimensional form for transmission as a sequential utterance, thereafter to be retranslated into some rough copy of the original in the head of the listener. Viewed in this light, it becomes entirely possible that there may exist no Type 1 theory of English syntax of the type that transformational grammar attempts to define -- that its constraints resemble wired-in conventions about useful ways of executing this tedious but vital operation, rather than deep principles about the nature of intelligence. An abstract theory of syntax may be an illusion, approximating what really happens only in the sense that Julesz's second order statistics theory approximates the behaviour of the set of processes that implement texture vision and which, in the final analysis, are all the theory that there is. In other words, the grammar of natural language may have a theory of Type 2 rather than of Type 1.

Even if a biological information processing problem has only a Type 2 theory, it may still be possible to infer more from a solution to it than the solution itself. This comes about because at some point in the implementation of a set of processes, design constraints attached to the machine in which they will run start to affect the structure of the implementation. This observation adds a different perspective to the two types of research carried out by linguists and by members of the artificial intelligence community. If the theory of syntax is really of Type 2, then any important implications about the CNS are likely to come from details of the way in which its constituent processes are implemented, and these are often explorable only by implementing them.

### **Conclusions**

If one accepts this view of A.I. research, one is led to judge its achievements according to rather clear criteria. What information processing problem has been isolated? Has a clean theory been developed for solving it, and if so how good are the arguments that

support it? If no clean theory has been given what is the evidence that favors a set-of-processes solution or suggests that no single clean theory exists for it, and how well does the proposed set of mechanisms work? For very advanced problems like story-understanding, current research is often purely exploratory. That is to say, in these areas our knowledge is so poor that we cannot even begin to formulate the appropriate questions, let alone solve them. It is important to realise that this is an inevitable phase of any human endeavor, personally risky (almost surely no exploring pioneer will himself succeed in finding a useful question), but a necessary precursor of eventual success.

Most of the history of A.I. (now fully 15 years old) has consisted of exploratory studies. Some of the best-known are Slagle's (1963) symbolic integration program, Weizenbaum's (1965) Eliza program, Evans' (1968) analogy program, Raphael's (1968) SIR, Quillian's (1968) semantic nets and Winograd's (1972) Shrdlu. All of these programs have (in retrospect) the property that they are either too simple to be interesting Type 1 theories, or very complex yet perform too poorly to be taken seriously as a Type 2 theory. Perhaps the only really successful Type 2 theory to emerge in the early phase of A.I. was Waltz's (1972) program. And yet many things have been learnt from these experiences -- mostly negative things (the first 20 obvious ideas about how intelligence might work are too simple or wrong) but including several positive things. The MACSYMA algebraic manipulation system is undeniably successful and useful, and it had its roots in programs like Slagle's. The mistakes made in the field lay not in having carried out such studies -- they formed an essential part of its development -- but consisted mainly in failures of judgement about their value, since it is now clear that few of the early studies even formulated any useful problems. Part of the reason for these internal failures of judgement lay in external pressures for early results from the field, but this is not the place to discuss what in the end are political matters.

Yet, I submit, one would err to judge these failures of judgement too harshly. They are merely the inevitable consequence of a necessary enthusiasm, based on a view of the long-term importance of the field that seems to me correct. All important fields of human endeavor start with a personal commitment based on faith rather than on results. A.I. is just one more example. Only a sour, crabbed and unadventurous spirit will hold it against us.

#### **Current trends**

Exploratory studies are important. Many people in the field expect that, deep in the heart of our understanding of intelligence, there will lie at least one and probably several important principles about how to organize and represent knowledge that in some sense captures what is important about the *general* nature of our intellectual abilities. An optimist might see the glimmer of such principles in programs like those of Sussman & Stallman (1975), of Marr & Nishihara (1975), in the overall attitudes to central problems set out by Minsky (1975), and possibly in some of Schank's (1973) work, although I sometimes feel that he almost missed the important points. While still somewhat cloudy, the ideas that seem to be emerging (and which owe much to the early exploratory studies) are:

- (1) That the "chunks" of reasoning, language, memory, and perception ought to be larger than

most recent theories in psychology have allowed (Minsky 1975). They must also be very flexible - at least as flexible as Marr & Nishihara's stick-figure 3-D models, and probably more. Straightforward mechanisms that are suggested by the terms "frame" and "terminal", are certainly too inflexible.

(2) That the perception of an event or of an object must include the simultaneous computation of several different descriptions of it, that capture diverse aspects of the use, purpose or circumstances of the event or object.

(3) That the various descriptions described in (2) include coarse versions as well as fine ones. These coarse descriptions are a vital link in choosing the appropriate overall scenarios demanded by (1), and in establishing correctly the roles played by the objects and actions that caused those scenarios to be chosen.

An example will help to make these points clear. If one reads

- (A) The fly buzzed irritatingly on the window-pane.
- (B) John picked up the newspaper.

the immediate inference is that John's intentions towards the fly are fundamentally malicious. If he had picked up the telephone, the inference would be less secure. It is generally agreed that an "insect-damaging" scenario is somehow deployed during the reading of these sentences, being suggested in its coarsest form by the fly buzzing irritatingly. Such a scenario will contain a reference to something that can squash an insect on a brittle surface -- a description which fits a newspaper but not a telephone. We must therefore conclude that when the newspaper is mentioned (or in the case of vision, seen) not only is it described internally as a newspaper, and some rough 3-D description of its shape and axes set up, but it is also described as a light, flexible object with area. Because sentence (B) might have continued "and sat down to read", the newspaper must also be described as reading-matter; similarly, as a combustible article, and so forth. Since one does not usually know in advance what aspect of an object or action is important, it follows that most of the time, a given object must be giving rise to several different coarse internal descriptions. Similarly for actions. It is important to note that the description of fly-swatting or reading or fire-lighting is not attached to the newspaper; merely that a description of the newspaper is available that will match its role in each scenario.

The important thing about Schank's "primitive actions" seems to me not the fact that there happens to be a certain small number of them, nor the idea that every act is expressed solely by reduction to them (which I cannot believe at all), nor even the idea that the scenarios to which they are attached contain all the answers for the present situation (that is where the missing flexibility comes in). The importance of a primitive, coarse catalogue of events and objects lies in the role such coarse descriptions play in the ultimate access and construction of perhaps exquisitely tailored specific scenarios, rather in the way



that a general 3-D animal model in Marr & Nishihara's theory can finish up as a very specific Cheshire Cat, after due interaction between the image and information stored in the primitive model. What after sentence (A) existed as little more than malicious intent towards the innocent fly becomes, with the additional information about the newspaper, a very specific case of fly-squashing.

Marr & Nishihara have labelled the problem of providing multiple-descriptions for the newspaper its "reference-window problem". Exactly how it is best done, and exactly what descriptions should accompany different words or perceived objects, is not yet known. These insights are the result of exploratory studies, and the problems to which they lead have yet to be precisely formulated, let alone satisfactorily solved. But it is now certain that some problems of this kind do exist and are important; and it seems likely that a fairly respectable theory of them will eventually emerge.

### *Mimicry versus exploration*

Finally, I would like to draw one more distinction that seems to be important when choosing a research problem, or when judging the value of completed work. The problem is that studies -- particularly of natural language understanding, problem-solving, or the structure of memory -- can easily degenerate into the writing of programs that do no more than mimic in an unenlightening way some small aspect of human performance. Weizenbaum (1976) now judges his program Eliza to belong to this category, and I have never seen any reason to disagree. More controversially, I would also criticise on the same grounds Newell and Simon's (1972) work on production systems, and some of Norman & Rumelhart's (1974) work on long term memory.

The reason is this. If one believes that the aim of information-processing studies is to formulate and understand particular information-processing problems, then it is the structure of those problems that is central, not the mechanisms through which they are implemented. Therefore, the first thing to do is to find problems that we can solve well, find out how to solve them, and examine our performance in the light of that understanding. The most fruitful source of such problems is operations that we perform well, fluently, reliably, (and hence unconsciously), since it is difficult to see how reliability could be achieved if there were no sound underlying method. On the other hand, problem-solving research has tended to concentrate on problems that we understand well intellectually but perform poorly on, like mental arithmetic; or on problems like crypt-arithmetic, geometry theorem-proving, or games like chess, in which human skills seem to rest on a huge base of knowledge and expertise. I argue that these are exceptionally good grounds for *not* studying how we carry out such tasks yet. I have no doubt that when we do mental arithmetic we are doing *something* well, but it is not arithmetic, and we seem far from understanding even one component of what that something is. Let us therefore concentrate on the simpler problems first, for there we have some hope of genuine advancement.

If one ignores this stricture, one is left in the end with unlikely looking mechanisms whose only recommendation is that they cannot do something we cannot do. Production systems seem to me to fit this description quite well. Anyone who has studied



them knows that as a programming language they are poorly designed, and hard to use. I cannot believe that the human brain could possibly be burdened with such poor implementation decisions at so basic a level.

A parallel may perhaps be drawn between production systems for students of problem-solving, and fourier analysis for visual neurophysiologists. Spatial frequency analysis of an image can mimic several interesting visual phenomena that seem to be exhibited by our visual systems. These include the detection of repetition, certain visual illusions, the notion of separate linearly adding channels, separation of overall shape from fine local detail, and a simple expression of size invariance. The reason why the spatial frequency domain is ignored by image analysts is that it is virtually useless for the main job of vision -- building up a description of what is there from the intensity array. The intuition that visual physiologists lack, and which is so important, is for how this may be done. Production systems exhibit several interesting ideas -- the absence of explicit subroutine calls, a blackboard communication channel, and some notion of a short term memory. But just because it displays these mechanisms (as Fourier analysis displays some visual illusions) does not mean that they have anything to do with what is really going on. My own guess would be, for example, that the fact that short-term memory can act as a storage register is probably the least important of its functions. I expect that there are several "intellectual reflexes" that operate on items held there, about which nothing is yet known, and which will eventually be held to be the crucial things about it because they perform central functions like opening up an item's reference window. Studying production systems seems to me a waste of time, because it amounts to studying a mechanism not a problem, and can therefore lead to no Type 1 results. The mechanisms that such research is trying to penetrate will be unravelled by studying problems, just as vision research is progressing because it is the *problem* of vision that is being attacked, not neural visual mechanisms.

A reflection of the same criticism can be made of Norman and Rumelhart's work, where they studied the way information seems to be organised in long term memory. Again, the danger is that questions are not asked in relation to a clear information-processing problem. Instead, they are asked and answers proposed in terms of a mechanism -- in this case, it is called an "active structural network" and it is so simple and general as to be devoid of theoretical substance. They may be able to say that such and such an "association" seems to exist, but they cannot say of what the association consists, nor that it has to be so because to solve problem X (which we can solve) you *need* a memory organized in such-and-such a way; and that if one has it, certain apparent "associations" occur as side-effects. Experimental psychology can do a valuable job in discovering facts that need explaining, including those about long-term memory, and the work of Shepard (1975) and of Warrington (1975) (for example) seems to me very successful at this; but like experimental neurophysiology, experimental psychology will not be able to explain those facts unless information-processing research has identified and solved the appropriate problems X. It seems to me that finding such problems X, and solving them, is what A.I. should be trying to do.

**Acknowledgement:** Although I take full responsibility for the purely personal views set out here, any virtues that they may have are due in part to many conversations with Drew McDermott.

### References

- Chomsky, A. N. (1965). *Aspects of the theory of syntax*. Cambridge, Mass.: M. I. T. Press.
- Cooley, J. M. & Tukey, J. W. (1965). An algorithm for the machine computation of complex Fourier series. *Math. Comp.*, 19, 297-301.
- Evans, T. (1968). A program for the solution of geomtric-analogy intelligence test questions. In *Semantic information processing*, Ed. M. Minsky, pp271-353. Cambridge, Mass.: M. I. T. Press.
- Horn, B. K. P. (1970). Shape from shading: a method for obtaining the shape of a smooth opaque object from one view. *Project MAC TR-79*.
- Jardine, N. & Sibson, R. (1971). *Mathematical taxonomy*. New York: Wiley.
- Julesz, B. (1975). Experiments in the visual perception of texture. *Scientific American*, 232, 34-43 (April issue).
- Levitt, M. & Warshel, A. (1975). Computer simulation of protein folding. *Nature*, 253, 694-698.
- Marr, D. (1975). Early processing of visual information. *M. I. T. A. I. Lab. Memo 340*.
- Marr, D. & Nishihara, H. K. (1975). Spatial disposition of axes in a generalized cylinder representation of objects that do not encompass the viewer. *M. I. T. A. I. Lab. Memo 341*.
- Minsky, M. (1975). A framework for representating knowledge. In: *The psychology of computer vision*, Ed. P. H. Winston, pp 211-277. New York: McGraw-Hill.
- Newell, A. & Simon, H. A. (1972). *Human problem solving*. Englewood Cliffs, N. J.: Prentice-Hall.
- Norman, D. A. & Rumelhart, D. E. (1974). *Explorations in cognition*. San Francisco: W. H. Freeman & Co. See e.g. the article entitled "The active structural network", pp 35-64.
- Quillian, M. R. (1968). Semantic memory. In: *Semantic information processing*, Ed. M. Minsky, pp227-270. Cambridge, Mass.: M. I. T. Press.

Raphael, B. (1968). SIR: semantic information retrieval. In: *Semantic information processing*, Ed. M. Minsky, pp33-145. Cambridge, Mass.: M. I. T. Press.

Schank, R. C. (1973). Identification of conceptualizations underlying natural language. In: *Computer models of thought and language*, Ed. R. C. Schank & K. M. Colby. San Francisco: W. H. Freeman.

Shepard, R. N. (1975). Form, formation, and transformation of internal representations. In: *Information processing and cognition: The Loyola Symposium*, Ed. R. Solso, pp 87-122. Hillsdale, N. J.: Lawrence Erlbaum Assoc.

Slagle, J. R. (1963). A heuristic program that solves symbolic integration problems in freshman calculus. In: *Computers and thought*, Ed. E. A. Feigenbaum & J. Feldman, pp191-203. New York: McGraw-Hill.

Sussman, G. J. & Stallman, R. M. (1975). Heuristic techniques in computer aided circuit analysis. *M. I. T. A. I. Lab. Memo 328*.

Ullman, S. (1975). On visual detection of light sources. *M. I. T. A. I. Lab. Memo 333*.

Waltz, D. L. (1972). Generating semantic descriptions from drawings of scenes with shadows. *M. I. T. A. I. Lab. Technical Report 271*.

Warrington, E. K. (1975). The selective impairment of semantic memory. *Quarterly Journal of Experimental Psychology*, 27, 635-657.

Weizenbaum, J. (1965). ELIZA - a computer program for the study of natural language communication between man and machine. *Communications of the Association for Computing Machinery*, 9, 36-45.

Weizenbaum, J. (1976). *Computer thought and human reason*. San Francisco: W. H. Freeman.